

English Abstract for JP 8-297677

Automatic method for generating thematic summaries - scoring each sentence of document based upon number of thematic terms contained within sentence, and selecting highest scoring sentences as thematic sentences

Patent Assignee: XEROX CORP (XERO)

Inventor: CHEN F R

Number of Countries: 005 Number of Patents: 003

Patent Family:

Patent No	Kind	Date	Applicat No	Kind	Date	Week
EP 737927	A2	19961016	EP 96302250	A	19960329	199646 B
JP 8297677	A	19961112	JP 9684297	A	19960405	199704
US 5689716	A	19971118	US 95422573	A	19950414	199801

Abstract (Basic): EP 737927 A

The method for automatically generating thematic summaries for machine readable representations of documents which include sentences and terms involves identifying thematic terms (42,44) within the document, by selecting as thematic terms a number of terms from the terms in the document, and scoring (46-58) each sentence of the document based upon the occurrence of thematic terms in each sentence. The highest scoring sentences are selected (62) as thematic sentences.

The selection of thematic terms involves determining a number of times each term occurs in the document, and selecting as thematic terms a number of terms from the terms in the document, based upon the number of times each term occurs in the document.

USE/ADVANTAGE - Automatically generating thematic summaries of documents. Automatically produces readable and semantically correct document summaries, and requires user to specify length of required summary. Document summaries may be automatically generated without using iterative approach.

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平8-297677

(43) 公開日 平成8年(1996)11月12日

(51) IntCl.

G 0 6 F 17/30

識別記号

庁内整理番号

9194-5L

F I

G 0 6 F 15/401

技術表示箇所

3 2 0 A

審査請求 未請求 請求項の数3 OL (全7頁)

(21) 出願番号 特願平8-84297

(22) 出願日 平成8年(1996)4月5日

(31) 優先権主張番号 4 2 2 5 7 3

(32) 優先日 1995年4月14日

(33) 優先権主張国 米国 (US)

(71) 出願人 590000798

ゼロックス コーポレーション

XEROX CORPORATION

アメリカ合衆国 ニューヨーク州 14644

ロチェスター ゼロックス スクエア

(番地なし)

(72) 発明者 フランシーヌ・アール・チェン

アメリカ合衆国 カリフォルニア州

94025 メンロパーク シャーマンアベニ

ュー 975

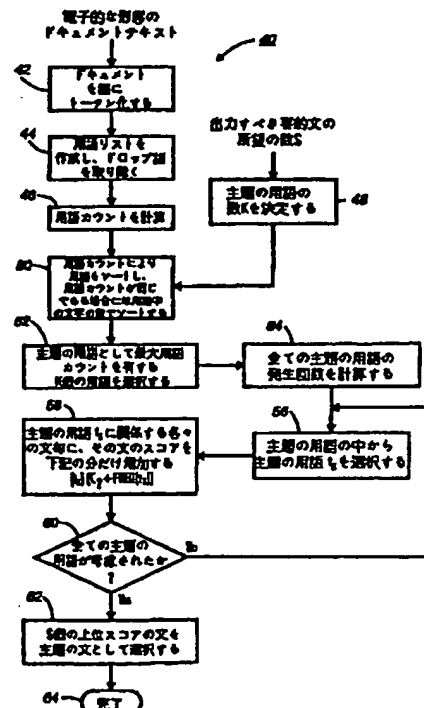
(74) 代理人 弁理士 小堀 益 (外1名)

(54) 【発明の名称】 主題の要約を生成する自動的な方法

(57) 【要約】

【課題】 機械で読み取り可能なドキュメントの主題の要約を自動的に生成する方法を提供する。

【解決手段】 第1の複数の文と第2の複数の用語を含んだドキュメントを機械で読み取り、プロセッサはプロセッサに結合されたメモリ内に電子的な形態で格納された命令を実行することにより、次のステップを実施する。a) 主題の用語として、前記第2の複数の用語から第1の数の用語を選択し、b) 各々の文の中の主題の用語の発生に基づいて前記第1の複数の文の各々の文にスコアを付け、c) 主題の文として、各々の文のスコアに基づいて前記第1の複数の文から第2の数の文を選択する。



【特許請求の範囲】

【請求項1】機械が読み取り可能な形態でプロセッサに対して提示されたドキュメントの主題の要約を生成するプロセッサにより実施される方法であって、ドキュメントは第1の複数の文と第2の複数の用語を含んでおり、前記プロセッサはプロセッサに結合されたメモリ内に電子的な形態で格納された命令を実行することにより前記方法を実施するものであり、

a) 主題の用語として、前記第2の複数の用語から第1の数の用語を選択するステップと、

b) 各々の文の中の主題の用語の発生に基づいて前記第1の複数の文の各々の文にスコアを付けるステップと、

c) 主題の文として、各々の文のスコアに基づいて前記第1の複数の文から第2の数の文を選択するステップとを含むプロセッサにより実施される方法。

—【請求項2】前記ドキュメントの中で主題の文が発生する順に主題の文をプロセッサのユーザに提示するステップを更に含む請求項1に記載のプロセッサにより実施される方法。

【請求項3】前記ステップb)が、文章の中で主題の用語が発生する度に、ドキュメントの中で主題の用語の発生頻度に関連した量だけ、各文のスコアを増加するステップを含む請求項1に記載のプロセッサにより実施される方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、自動的なテキスト処理の方法に関する。特に本発明は、ドキュメントの主題の要約を生成する自動的な方法に関する。

【0002】

【従来の技術】ドキュメント要約及び概要は、ドキュメントを検討するのに必要な時間を減少させることによって有効な機能を果たす。要約及び概要は、ドキュメント作成の後に手動或いは自動的に生成することができる。手動の要約及び概要は、高品質であるが人間の労働が必要であるので高価になるおそれがある。別法として、要約及び概要は、自動的に生成することができる。自動的な要約及び概要は、安価に作成することができるが、高品質を一貫して得ることは困難である。

【0003】

【発明が解決しようとする課題】自動的な要約を生成するためのシステムは、二つの計算機的な技法、すなわち、自然言語処理、或いは、量的な内容分析の二つに頼っている。自然言語処理は、計算機的な処理を集中的に行う。これに加えて、ドキュメント内容が限定されていないときには、自然言語処理を使って意味的に正しい要約及び概要を作成することは困難である。

【0004】量的な内容分析は、テキストの統計上の特性に頼って要約を作成する。ジェラルド・サルトン(Gerald Salton)は、「自動テキスト処理(Automatic Text P

rocessing)」(1989)において、ドキュメントを要約するための量的な内容分析の使用について議論している。サルトン要約器(Salton summarizer)は、最初に、ドキュメント集成の中のテキスト語を分離する。次に、サルトン要約器は、タイトル、図、キャプション、脚注において使用された語を、タイトル語としてフラグを立てる。その後、ドキュメント集成の中の残りのテキスト語の発生頻度が決定される。次いで、発生頻度とテキスト語の位置は、語の重みを生成するために使われる。サルトン要約器は、語の重みを使用して、ドキュメント集成の中の各々のドキュメントの各々の文にスコアを付ける。これらの文スコアは順番に使用されて、ドキュメント集成の中の各々のドキュメント毎に、所定の長さの要約を作成する。語の重みは、各々の個別のドキュメントの中よりは、ドキュメント集成の全体での発生に基づいて決定されるので、個別のドキュメントのテーマを正確には反映していない恐れがある。

【0005】

【課題を解決するための手段】機械で読み取り可能なドキュメントの主題の要約を自動的に生成する技法が説明される。この技法は、ドキュメント内の主題の用語の識別で始まる。次に、ドキュメントの各々の文に、その文に含まれる主題の用語の数に基づいて、スコアが付けられる。その後、最もスコアが高い文が、主題の文として選択される。

【0006】

【発明の実施の形態】図1は、自動的にドキュメントの主題の要約を生成するためのコンピュータシステムを示す。

【0007】図2は、図1のコンピュータシステムを使用するドキュメントの主題の要約を生成する方法のフローチャートである。

【0008】図1は、本方法が実施されたコンピュータシステム10を、ブロック図形式で示す。本方法は、コンピュータシステム10の動作を変え、機械が読み取り可能な形態で表されたどのようなドキュメントの主題の要約も生成することを可能にする。簡単に説明されたように、コンピュータシステム10は、ドキュメント内で主題の用語を識別し、次いで、文内に含まれる主題の用語の数に基づいてドキュメントの各々の文にスコアを付けることにより、主題の要約を生成する。その後、コンピュータシステム10は、最も高いスコアが付けられた文を、主題の文として選択し、そしてそれらの文をコンピュータシステム10のユーザに提示する。

【0009】本方法のより詳細な議論の前に、コンピュータシステム10について検討する。コンピュータシステム10は、コンピュータユーザへ情報を視覚的に表示するためのモニタ12を含む。コンピュータシステム10は、プリンタ13によってもコンピュータユーザに情報を出力する。コンピュータシステム10は、コンピュ

ータユーザに、データを入力するための幾つもの方法を提供する。キーボード14は、コンピュータユーザがタイピングによってコンピュータシステム10にデータを入力することを可能にする。マウス16を移動することによって、コンピュータユーザは、モニタ12に表示されたポインタを移動することが可能になる。また、コンピュータユーザは、スタイラス或いはペンで電子タブレット18に書くことによって、コンピュータシステム10に情報を入力することができる。別法として、コンピュータユーザは、フロッピーディスクドライブ22にディスクをさしこむことによって、フロッピーディスクのような磁気媒体に格納されたデータを入力することができる。光学的文字認識ユニット(OCRユニット)24により、コンピュータユーザがハードコピードキュメントをコンピュータシステムに入力することが可能になり、次いで、OCRユニット24は、符号化された電子的表現、典型的には情報交換用米国標準コード(ASCII)に変換される。

【0010】コンピュータユーザの命令を実行するために、プロセッサ11は、コンピュータシステム10の動作を制御すると共に調整する。メモリに電子的に、すなわち、メモリ28或いはディスクドライブ22内のフロッピーディスクのいずれかに、格納された命令を実行することによって、プロセッサ11は、各々のユーザコマンドに応じた適切な挙動を決定し且つ行う。典型的には、プロセッサ11に対する動作命令は、固体メモリ28に格納され、命令への頻繁で迅速なアクセスが可能となる。使うことができる半導体メモリ装置には、読み出し専用メモリ(ROM)、ランダムアクセスメモリ(RAM)、ダイナミックランダムアクセスメモリ(DRAM)、プログラム可能な読み出し専用メモリ(PROM)、消去可能なプログラム可能な読み出し専用メモリ(EPROM)、フラッシュメモリのような電氣的に消去可能なプログラム可能な読み出し専用メモリ(EEPROM)が含まれる。

【0011】図2は、機械が読み取り可能な主題の要約を生成するために、プロセッサ11によって実行された命令40をフローチャート形式で示す。命令40は、固体メモリ28内に、或いは、フロッピーディスクドライブ22内に置かれたフロッピーディスクに格納することができる。命令40は、LISPとC++を含むどのようなコンピュータ言語でも実現することができる。

【0012】命令40の実行を開始するためには、ドキュメントを電子的な形態で選択して入力することが必要である。もし所望であるならば、命令40の実行開始前に、コンピュータユーザは、「S」で示された主題の要約の長さを、デフォルトの長さから変えてもよい。主題の要約のデフォルトの長さは、任意の数の文に設定することができる。ドキュメントの拾い読みを意図している実施態様においては、主題の要約のデフォルトの長さ

は、五つの文に設定される。

【0013】プロセッサ11は、ステップ42に分岐することによって、要約すべきドキュメントの選択にตอบสนองする。ステップ42の期間中は、プロセッサ11は、選択された語及び文をトークン化する。すなわち、プロセッサ11は、選択されたドキュメントの、機械が読み取り可能な表現を分析し、文の境界及び各文の中の語を識別する。自然言語テキストのトークン化は周知であり、したがって、ここでは詳細には説明されない。これに加えて、トークン化の期間中は、プロセッサ11は、ドキュメントの各々の文に文I. D. を割り当てる。一つの実施態様においては、各々の文は、ドキュメントの開始に関してその位置を表している数によって識別される。文を識別する他の方法が、本方法に影響を与えることなく使用できる。選択されたドキュメントのトークン化の後で、プロセッサ11はステップ42からステップ44へ分岐する。

【0014】プロセッサ11は、ステップ44の間に、ドキュメントの各々の語トークンを調べ、その語を用語リストに既に含まれている用語を比較する。語トークンがリストにまだ含まれていない場合には、次いで、プロセッサ11は、その語を用語リストに加えて、その語が発生した文の文I. D. を注記する。他方、語が用語リストに既にある場合には、プロセッサ11は、その用語についてのエントリ或いはリストへ、その語についての文I. D. を単純に加える。言い換えれば、ステップ44の期間中、プロセッサ11は、その言葉の発生毎の位置とドキュメントの語を関連させるデータ構造を生成する。このように、たとえば、「背教(apostasy)、7、9、12」の用語リストエントリは、用語「背教」が、ドキュメントの文7と9と12で発生するというを示す。

【0015】好ましくは、用語リストを生み出している間、プロセッサ11は、ストップ語を取り除く。ここで使用されたように、「ストップ語(stop word)」は、主題の意味を伝達せず、また、自然語テキストにおいて非常に頻繁に発生する語である。ほとんどの代名詞、前置詞、決定詞、及び、動詞「である(to be)」は、ストップ語として分類される。このように、例えば、「そして(and)、一つの(a)、その(the)、～の上の(on)、～によって(by)、～について(about)、彼(he)、彼女(she)」のような語は、ストップ語である。ドキュメント内のストップ語は、ストップ語のリストとドキュメントについての語トークンを比較することによって識別される。用語リストからストップ語を削除することは必要でないが、削除すれば、ドキュメントの主題の要約を生成するのに必要な全体の処理時間が減少する。

【0016】プロセッサ11は、用語リストを完成した後でステップ44からステップ46へ分岐する。ステップ46の間に、プロセッサ11は用語リストを分析し

5

て、ドキュメント中で各用語の発生回数を決定する。これは、単純に、その言葉に関係している文I. D. の数を数えることによって行なわれる。それが行なわれて、プロセッサ11は、ステップ50に分岐する。

【0017】実行の開始より後で、ステップ50の実行の前に、ステップ48の間に、プロセッサ11は、主題の文を選ぶ際に使用されるべき主題の用語の数を決定する。「K」で示されたその数は、主題の要約の長さに基づいて、すなわち、Sに基づいて、決定される。一般に、KはS未満で1より大きくあるべきである。KをSより小さくすることにより、選択された主題の文の間の幾分かの共通性を確実にする。好ましくは、Kは下式に従って決定される。

【0018】

【数1】

$$K = \begin{cases} S \times c_1 & 8 \times c_1 > 3 \\ 3 & \text{その他の場合} \end{cases}$$

ここで、 c_1 はその値が1未満の定数、Sは主題の要約の中の文の数、Kは主題の用語の数である。

【0019】一つの実施態様においては、 c_1 の値が、0.7と等しくされる。

【0020】Kの値とステップ46の間に生成された用語カウントが与えられ、プロセッサ11は、K個の主題の用語を選択する処理を始める。ステップ50の間、プロセッサ11は、それらのカウント、すなわち、ドキュメントの中の各々の用語の総発生回数に従って用語リストの用語をソートする。二つの用語が同じカウントを有する場合には、最大文字数を含む用語の方が選択される。ソートされた用語リストが生成され、そのリストがメモリに格納されると、プロセッサ11は、ステップ50からステップ52に分岐する。ステップ52の間に、プロセッサは、ソートされた用語リストから最も高いカウントを有するK個の用語を選択する。それが行なわれて、プロセッサ11は、ステップ54に進む。

【0021】ステップ54の間に、プロセッサ11は、ドキュメントの中のK個の主題の用語の総発生回数を計算する。「N」で示されるその数は、K個の主題の用語のカウントを合計することによって計算される。プロセッサ11は、ステップ54からステップ56に分岐する。

【0022】主題の用語が選択され、それらのカウントが決定されると、プロセッサ11は、ドキュメントの文の主題の内容を評価することを始める準備ができる。ステップ56、58、60、62の間に、プロセッサ11は、K個の主題の用語の少なくとも一つを含むそれらの文だけを考慮する。プロセッサ11は、ソートされた用語リストのK個の最も高いスコアが付けられた用語を調べることによってそれを行う。 t_s で示された用語を選択した後に、ステップ56の間に、プロセッサ11は、

6

ステップ58の間に t_s に関係している各々の文I.

D. を調べる。 t_s に関係している各々の文I. D. について、プロセッサ11は、文のスコアを増加する。好ましくは、各々の文毎のスコアが、下式で表されるsだけ増加される。

【0023】

$s = \text{count } t_s [c_2 + \text{freq } t_s];$

ここで、 $\text{count } t_s$ は、文の中の t_s の発生回数、 c_2 は、ゼロでない正の値を有する定数、 $\text{freq } t_s$ は、選択された用語 t_s の頻度である。 $\text{freq } t_s$ は、下式で表される。

$\text{freq } t_s = \text{count } t_s / N;$

ここで、Nは、ドキュメント内の主題の用語の総発生回数を表す。好ましくは、 c_2 は1の値に設定される。

【0024】文スコアは、ステップ58の間に文スコアリストを生成することによって跡を追うことができる。プロセッサ11が文I. D. を選択する毎に、文スコアリストが調べられ、その文I. D. を含んでいるかどうかを見る。含んでいない場合には、その文I. D. が文スコアリストに追加され、そのスコアが適切に増加される。他方、文スコアリストが既に特定の文I. D. を含んでいる場合には、次いで、その文と既に関係があるスコアが、先に述べたような方法で増加される。

【0025】選択された用語 t_s に関係する全てのスコアを増加した後に、プロセッサ11は、ステップ58からステップ60へ分岐する。ステップ60の間、プロセッサ11は、全ての主題の用語が評価されたかどうか決定する。そうでない場合には、プロセッサ11はステップ56へ戻って、選択された用語として他の主題の用語を選択する。プロセッサ11は、すべての主題の用語が調べられるまで、先に説明したように、ステップ56、58、60を通して分岐する。その事象が発生するとき、プロセッサ11は、ステップ60からステップ62に分岐する。

【0026】ステップ62の間、プロセッサ11は、主題の要約として、最も高いスコアを有するS個の文を選択する。プロセッサ11は、スコアによって文スコアリストをソートすることによってこれを行なう。主題の文が選択されると、プロセッサ11は、主題の要約をユーザにモニタ12或いはプリンタ13を介して提示することができ、また、主題の要約を後で使用するためにメモリやフロッピディスクに格納することもできる。主題の要約の文は、好ましくは、ドキュメント内の発生順に提示される。文は、パラグラフ形式で提示することができるが、文は論理的にパラグラフを形成しないかもしれないので、各々の文について個別に提示することが好ましい。主題の要約の生成が完了すると、プロセッサ11は、ステップ62からステップ64に分岐する。

【0027】このように、ドキュメントのための主題の要約を自動的に生成する方法が説明された。この方法

50

は、量的な内容分析に頼って主題の用語を識別し、これは次に主題の文を識別するために使用される。付録Aと付録Bは、主題の要約を自動的に生成する方法を使って生成された要約を含む。

【0028】付録A：シュワルナゼの辞任演説の要約
私がこのような演説のテキストを作成し、私がそれを書記局に渡し、そして代理人がそれを知ることができる一
国の指導者によって、大統領によって、そして外務省
によって行われたきた現在の政策の範囲は何か、また、
国の発展、我々の民主化と国の再生、経済発展等のため
に、どのようにして現在の状況が形作られてきたかであ
る。

【0029】昨日、何人かの同志による演説があった。
――彼らは退役軍人である。――彼らは、大統領と国の
指導者がベルシャ湾へ軍隊を送ること禁止することを採
択する宣言の必要性に疑義を呈した。そして、これらの
昨日の同志の演説は、忍耐の杯を満たし、そして溢れさ
せた。

【0030】国内及び海外の双方で、10回程、私はこの
対立に対するソビエト連邦の態度を講演し、説明しな
ければならなかった。

【0031】その場合には、新しい政治的な思想の原則
を主張する分野において、我々の全てによって、国全体
によって、及び、我々の人民の全てによって、近年行わ
れた全てのことに衝突しなければならなかった。

【0032】第2に、私がくり返して説明し、そして、
ミハイル・セルゲイビッチが最高議会における彼の演説
でこれを話したように、ソビエト指導者は何も計画を持
っていない。――少なくとも私は知らない。多分誰かが
何らかの計画を持っているかもしれないが、あるグルー
プ――但し公の機関、国防省――外務大臣がベルシャ湾
付近で軍隊を上陸させる計画を立てることが非難され
た。

【0033】第3の問題は、私がそこで言ったことであ
り、そして、私が確認し公に述べたことであるが、ソビ
エトの人民の利益が侵略される場合には、ただ一人の人民
が被害を受ける場合であっても、どこで起きた場合であ
っても、どの国においてであっても、イラクにおいて
だけではなく、他のどの国においても、――無論、ソビ
エト政府においても、ソビエト側の意思は、その市民の
利益を擁護する。

【0034】それでもやはり、私は、これは偶発的な事
象ではないと主張する。失礼、私はいまソビエト最高会
議の議会を召集するところである。同志ルカノフのイニ
シアティブで、文字通り会議の前に、ドイツ民主共和国
との条約に関する協議事項に重大な問題が含まれた。

【0035】私は、私の国で起こっていること、そし

て、我々の人民を待つ裁判に甘んじることはできない。

【0036】付録B：ジョン・シーリー・ブラウン(John Seely Brown)による「会社を再発明する研究(Research that Reinvents Corporation)」の要約

会社が技術の急速な変化にペースを合わせ、そして、不安定なビジネス環境に対処しようとするとき、研究部門は、単に新製品を発明すること以上のことを行わなければならない。

【0037】次の十年間で、PARCの研究員は、パーソナルコンピュータの革命的であるだけの基本的な発明のいくつかに対して責任があり、他の会社がゼロックスより速くこれらの発明を商業化するのを見ていた。

【0038】これらの問題に対する一つの人気がある解答は、研究部門の焦点を革新的なブレイクスルーから離れて、段階的な発明の方向に移動させることであり、基本的な研究から離れて応用研究の方向に移動させることである。

【0039】パイオニア的な研究を行うことを我々が強調することが、我々に技術、発明、そして実際研究自身が何を意味するのかを再定義させた。

【0040】そのような活動は、情報技術「遍在する計算(ubiquitous computing)」、すなわち、広範囲な毎日の目的における情報技術の組み込みにおける次の偉大なブレイクスルーを成功裏に開発するために、会社にとって不可欠である。

【0041】会社組織の研究が、その製品だけでなく会社の営業に集中することを始めるとき、他の原則は急速に明確になる。発明は、研究部門の特権的な活動ではない。PARCにおいて、我々は、ゼロックスのビジネスの最前線の従業員による局地的な発明のこの処理を検討し、技術を発展させて、全体として会社のための収穫を得るための技術を開発する。

【0042】結果：ゼロックスのコア製品への重要な貢献だけでなく、我々の会社を遙に越えて実施される発明への独特なアプローチ。

【図面の簡単な説明】

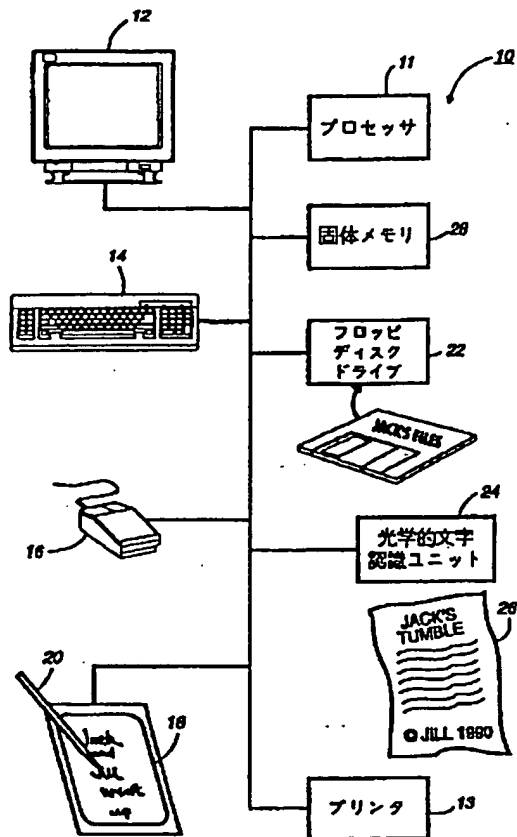
【図1】 自動的にドキュメントの主題の要約を生成するためのコンピュータシステムを示す。

【図2】 図1のコンピュータシステムを使用するドキュメントの主題の要約を生成する方法のフローチャートである。

【符号の説明】

10 コンピュータシステム、11 プロセッサ、12 モニタ、13 プリンタ、14 キーボード、16 マウス、18 タブレット、22 ディスクドライブ、24 OCRユニット、28 固体メモリ

【図1】



【図2】

